

TEI and the Nxaʔamxcín Dictionary

Ewa Czaykowska-Higgins and Martin Holmes
University of Victoria

Acknowledgements

- The late M. Dale Kinkade
- The speakers who worked with Dale in the 1960s & 70s
- Mary Marchand, Elizabeth Davis, the late Agatha Bart and the late Tillie George who urged us to work on a dictionary in 1990
- Mary Marchand and Pauline Stensgar who continue to support this project
- The various students who worked on inputting the materials into Lexware and XML, and programmers Bob Hsu and Greg Newton
- Colville Tribes Culture Committee
- Social Sciences and Humanities Research Council of Canada, University of Victoria Faculty of Humanities, and Humanities Computing and Media Centre

Outline of Presentation

Part I Background: The Project and The Nxaʔamxcín Language

Early Beginnings of the Project: From Filecards to Lexware

Part II Using TEI for a Lexical Resource Project

Evolution and Early Challenges

From Lexware to Text-Encoding Initiative structure

Describing what TEI can do

Part III Comparing Standards

Format and Terminological Interoperability



Part I

Background The Project and The Nxaʔamxcín
Language

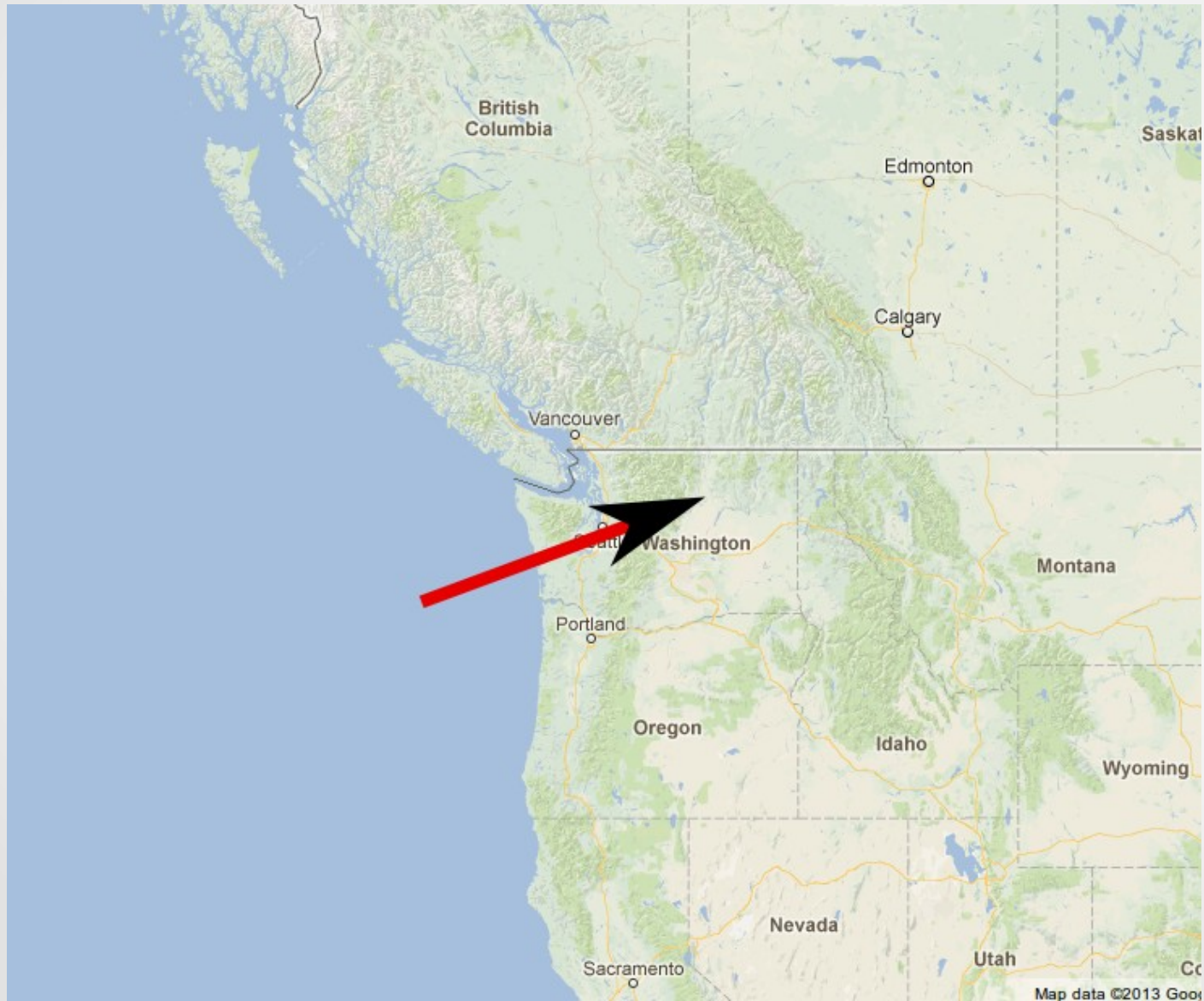
The Project Today

- To construct and make publicly available
 - Rich and extensive lexical documentation of Nxaʔamxcín
 - Recorded in the 1960s and 70s with approx 22 speakers, some in their 80s and 90s at the time
 - Therefore: legacy materials
- To create from these legacy materials lexical resources
 - That are both enduring and reusable in various formats and for various purposes
 - Including as a print dictionary and as a searchable, extendable web-based database which can be used in different views, for learning, teaching, language revitalization, and scholarly purposes

The Project Team

- At Colville Tribes
 - Elders: Pauline Stensgar, Mary Marchand
 - Language Program: Ernest K'saw's Brooks, Sharon Covington; Albert Andrews
 - History and Archeology: Guy Moura
- At UVIC
 - Martin Holmes, programmer
 - Sarah Kell, editor
 - Caitlin Bird, research assistant
 - Ewa Czaykowska, compiler and chief editor

Map



The Nxaʔamxcín Language

- Southern Interior Salish; one of 7 Interior languages; one of 23 languages in the Salish family
- Spoken historically in north central Washington State
- Since 1870, confined to the Colville Reservation by Presidential Executive Order together with 8 other nomadic tribes
- Four varieties: Wenatchi, Moses-Columbia, Entiat and Chelan; roughly 3500 tribal members descended from Nxaʔamxcín-speaking families
- Fluent and active speakers in 2013: 2
- Nxaʔamxcín Language Program-since late 1990s

Initial purpose of project: 1991

- Based on Kinkade's fieldnotes and how he had transformed these notes into “data” on filecards
- To create a print dictionary
- Containing rich morphological information
- Organized according to root morphemes and lexicalized stems

The initial target dictionary

- The entries were intended to look like the following from Kinkade's Upper Chehalis dictionary (1991)

628. $\sqrt{k^w}\text{əná-}$: hold, take, take hold, get, pick up, grab.
1. $\sqrt{k^w}\text{aná-m'ł}$, $s\sqrt{k^w}\text{aná-m'al-n}$ *detr.* 2. FB
 $s\sqrt{k^w}\text{aná?-m'ł}$, -- *detr. pl.* 3a. $\sqrt{k^w}\text{aná-n}$,
 $s\sqrt{k^w}\text{aná-t-n}$ *tr.* ME --, $kwâ-nât$ *tr.* 3b. TC
 $k^w\text{anátn}$ grasped, got hold. 3c. TC $łk^w\text{anáčs}$ *refl.*
you'll take. 4. FB $\sqrt{k^w}\text{an'á-n}$ or $\sqrt{k^w}\text{aná?-n}$,
 $s\sqrt{k^w}\text{an[a?]}á-t-n$ or $s\sqrt{k^w}\text{an[?]}á-t-n$ *tr. pl.* 5. TCh FB
 $?ac\sqrt{k^w}\text{ən-ł}$ he held. 6. ME $kwĩn-tûk$ hold. 7. FB $?it$
 $\sqrt{k^w}\text{aná?-n}$ he held it for a while. 8. $\sqrt{k^w}\text{əna-}$: catch,
catch hold of, get, detain, FB capture. 8a. $\sqrt{k^w}\text{əna-x}^w$ or
FB $\sqrt{k^w}\text{əná-x}^w$, $s\sqrt{k^w}\text{əna-y-n}$ *tr.* 8b. FB $\sqrt{k^w}\text{ən'á-x}^w$,
 $s\sqrt{k^w}\text{ən'á-y-n}$ *tr. pl.* 9. FB $\sqrt{k^w}\text{ána?-ł-čł}$ she had,

The filecards



Filecard closeup

Cm

sén quiet person Y
 gentle Y24.42

sénsènt tame, gentle Y17.23; Y18.126; Y28.42; W10.85
 quiet person Y24.40

sénp tame, get tame Y28.43; EP
 he got gentle W10.86

sénsèntwìx he got gentle W10.87

tēl sénsènt he's real gentle W10.90

tīl sēnp he's gentle now W10.91

sēnpmūnən I tamed him W10.169

sēnpstūnən I tamed him W10.170

sēsənt tame EP2.68.8 / it's tame/gentle JM3.20.11

sēsəntləx they're tame/gentle JM3.21.1

cf. qémqèmt quiet person

ca. sañ-t

Filecards to Lexware

- In 1991
- The “data” were input into a computer program called Lexware
- Lexware developed by Robert Hsu at U of Hawaii to create lexical resources
- Used band format, including bandnames, headwords, numbering and spaces to create a hierarchical system of embedding sub-entries within entries (Hsu 1985)

Lexware entry

.rt √sən
g *quiet person, *gentle
k Y; Y24.42

..ch
| infl stative
| (stt √sən+sən-t
| | g *tame, *gentle
| | k Y17.23; Y18.126; Y28.42;
| | var √sən+sən-t
| | g quiet person
| | k Y24.40
| | var √sən+sən-t
| | g *tame
| | k EP2.68.8
| | var √sən+sən-t
| | g it is *tamē or *gentle
| | k JM3.21.11

| 2stt √sən+sən-t ləx
| 2g they are *tame or *gentle
| 2k JM3.20.11

| 2il.ch tət √sən+sən-t
| 2df he is real *gentle
| 2k W10.90

..in √sən-p
g *tame, get tame
k Y28.43; EP
var √sən-p
g he got gentle
k W10.86

il.in t'íl' √sən-p
df he is *gentle now
k W10.91

infl success
nn √sən-p-nún-ən
g I tamed <*tame> him
k W10.169
q MDK underlined root schwa,

The data

.rt? ÀEOTBELÀÀ1SOHÀkÀSOEOTÀ

linfl transitive

1ltr ÀEOTBELÀÀ1SOHÀÀUSOHÀkÀSOEOTÀ©À9SOHÀn

1lg *pull

q should this be under a separate root?

The rescue

- Lexware binary data was dependent on:
 - WordPerfect for DOS
 - Customized character sets
 - Obsolete printer fonts
 - A Hercules graphics card
- AARGH! 3 months of misery. See Holmes & Newton (2008) for the whole horrible story.

Lessons learned

- Hardware goes away.
- Software goes away.
- Only your data remains.
- You'd better be able to read it.



Part II

Using TEI for a Lexical Resource Project

Never again...

- From now on:
 - **Unicode** (stable, mature standard)
 - **XML** (stable, mature standard with conversion tools)
 - **TEI**

TEI: Text Encoding Initiative

- “...a consortium which collectively develops and maintains a standard for the representation of texts in digital form”
- Chief deliverables: schemas and a set of *Guidelines* which specify encoding methods for machine-readable texts
- Regular releases, thousands of projects, lots of support
- **Chapter 9: Dictionaries:**
 - Tags and attributes for encoding historical, print and born-digital lexica

TEI

543 elements,
468 attributes

DTD

RNG

RNC

XSD

PDF

HTML

ePub

Documentation

Schemas

<?xml...

ODD file:
customize,
constrain,
document

<person xml:id="MDH" role="">

<persName>

<surname>Holmes</surname>

</persName>

<affiliation>University of Victoria</affiliation>

</person>

</listPerson>

articDesc>

fileDesc>

▼ Editor

▼ Elder

▼ FieldResearcher

▼ FluentSpeaker

▼ LangApprentice

▼ LangTeacher

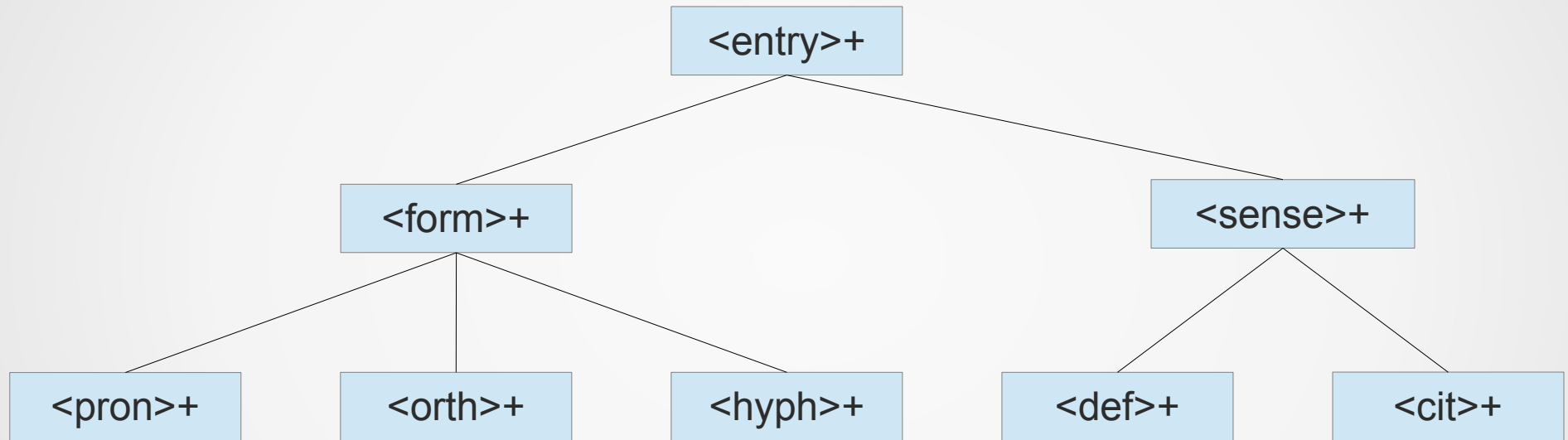
▼ ProgramManager

▼ Programmer

(Computer Programmer on the project)

g and Media Centre</affiliation>

Structure of a TEI entry



Morpheme linking

- All morphemes have entries with feature structures:

```
<entry xml:id="ḥam?">
  <form>
    <pron>
      <seg type="p" subtype="i">ḥám?</seg>
      <bibl corresp="psn:ECH">ECH</bibl>
    </pron>
  </form>
  <sense>
    <def>
      <seg><gloss>forbid</gloss>, <gloss>stop</gloss></seg>
      <note type="editorial" resp="psn:ECH">Definition inferred based on complex forms containing this
root.</note>
    </def>
  </sense>
  <fs>
    <f name="baseType">
      <symbol value="root"/>
    </f>
  </fs>
  <note type="editorial" resp="psn:ECH">Root entry added based on attested complex forms.</note>
</entry>
```

Morpheme linking

- All polymorphemic items are linked to the entries of their constituent morphemes:

```
<entry xml:id="hámʔcn">
  <form>
    <pron>
      <seg type="p" subtype="i">hámʔcn</seg>
      <bibl corresp="psn:ECH">ECH</bibl>
      <seg type="n">hámʔčən</seg>
      <bibl corresp="psn:JM">JM3.200.10</bibl>
    </pron>
    <hyph>
      √<m corresp="m:hámʔ">hámʔ</m>-<m corresp="m:stu">c</m>-<m corresp="m:Ø-OBJ m:n-
SUBJ">n</m>
    </hyph>
  </form>
  [...]
</entry>
```


So we can do this:

/hámʔcn/ ECH [hámʔčən] JM3.200.10

Get related words

√hámʔ -c -n

COMPONENT MORPHEMES:

hámʔ

c

Ø-OBJ

n-SUBJ

/hámʔ/ (ROOT) ECH

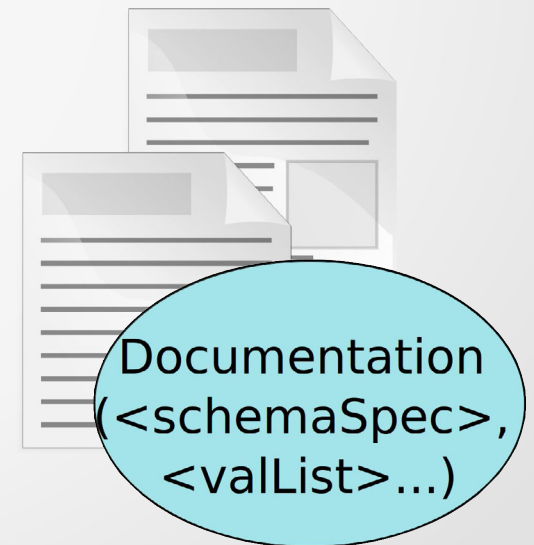
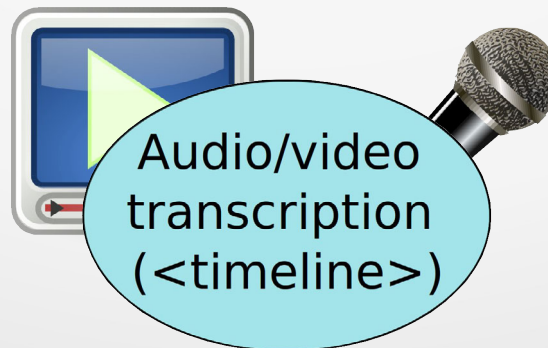
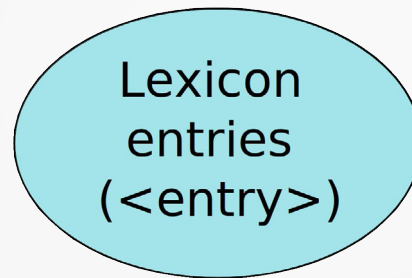
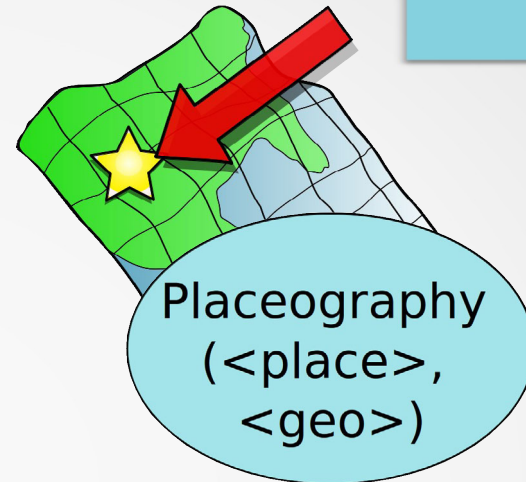
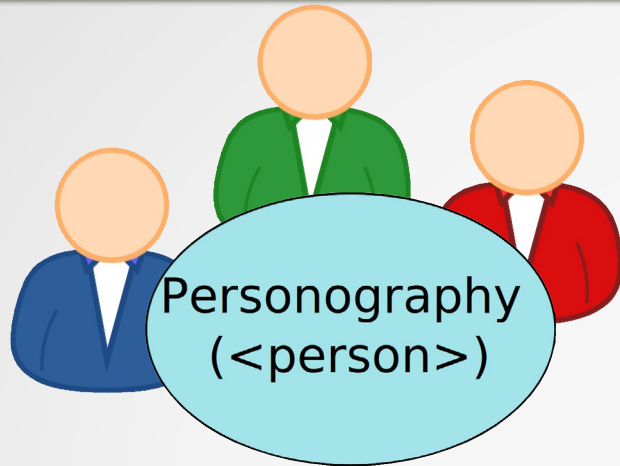
forbid, stop Note: Definition inferred based on complex forms containing this root.

Other entries containing this morpheme

| | | |
|---|------------|-----------------------|
| ■ | ʔachámʔsn | stopped |
| ■ | hámʔcn | stopped |
| ■ | hámʔn | stop, forbid, forbade |
| ■ | hámʔn lx | stopped |
| ■ | hámʔntm lx | stopped |

I stopped him from arguing/quarreling with someone JM3.200.10

TEI





Part III

Comparing Standards

Standardization and Interoperability

- Bird and Simons (2003) “Seven Kinds of Portability”

“If digital language documentation and description should transcend time, they should be reusable... across different software and hardware platforms....and across different purposes.”

- Interoperability is the ease of moving between systems/platforms
- We use it here to refer to ease of moving between TEI and other dictionary-creating formats
- Interoperability requires standardization of format and terminology

Format Interoperability

- XML-to-XML conversion is trivial.
- What schema should we aim at?
- **LMF version 16 DTD?**
 - 5 drawbacks at <http://tla.mpi.nl/relish/lmf/>.
- **LIFT?**
 - Supported by *WeSay*, *FLEx* and *Lexique Pro*, but last spec 2009, last code update over a year ago.
- **RELISH?**
 - Best candidate... Modular RNG schemas, and allows for TEI feature structures.

RELISH encoding

```
<LexicalEntry>
  <tei:f name="partOfSpeech"
    dcr:datcat="http://www.isocat.org/datcat/DC-1345"
    fVal="commonNoun"
    dcr:valueDatcat="http://www.isocat.org/datcat/DC-1256"/>
  <Lemma type="Form">
    <tei:f name="writtenForm"
      dcr:datcat="http://www.isocat.org/datcat/DC-1836"
      fVal="clergyman"/>
  </Lemma>
</LexicalEntry>
```

Terminological Interoperability

- RELISH depends on feature structures
- Feature structures depend on a shared ontology
- And that's a complication...
- One possible ontology is the GOLD ontology (General Ontology for Linguistic Description; <http://linguistics-ontology.org>), an attempt to create a stable, structured and documented standard
- GOLD is available through the ISOcat Data Category Registry which provides a formal system for associated terms from ontologies with feature structures in a specific encoding, using ISOcat URIs associated with specific concepts.

GOLD, ISOcat and Nxaʔamxcín

- Our feature structures are based on a morphological analysis of Nxaʔamxcín that itself is based on
 - 1) traditional Salishanist analysis as represented by the work of Kinkade;
 - 2) modified by subsequent analyses found in Willett (2003) and Czaykowska-Higgins (1998)
- This morphological analysis proposes a particular set of categories and relationships between the categories

(Mis)matches between ontologies

- We can match a small subset of our features straightforwardly onto the GOLD ontology
- There are also cases which have no match in GOLD but can be added
- But, there is a large set of cases (most) which map onto GOLD only loosely; crucially these mismatches arise from our specific morphological analysis
 - e.g., Distributive in Nxaʔamxcín involves plurality and events/actions repeated or distributed over time/space; the closest match is AugmentativeSize, but Distributive is not a Size category in Nxaʔamxcín

(Mis)matches between ontologies

- Nxa'amxcín words have different suffixes that signal that stems are transitive or intransitive (e.g., applicatives); these are categorized in the Salishanist literature as 'Valence-changing' morphemes; GOLD categorizes many of these as Voice.
- The supracategories in GOLD do not align with the supracategories in our ontology.
 - e.g., GOLD's Voice properties (like applicative) are included in a supracategory 'Morphosyntactic'; its Transitivity category is included in a supracategory 'Derivational'. But in the Nxa'amxcín ontology all transitivity markers are in one supracategory 'Morphosyntactic:Valence'
 - While it is not impossible to create a shared ontology, it is also not straightforward, and it is time-consuming

Conclusion

- Ultimately Interoperability is an important long term goal.
- Format conversion (XML → XML) is trivial, but...
- True interoperability depends on a shared ontology.
- The shared ontology depends on a shared grammatical analysis.
- This will take time to develop.
- In the meantime, in TEI, we have a stable, flexible schema system with solid documentation and a good toolset for creating the lexical resources we need now.



Lámlamt

To Unicode...

Transformer Unicode Replace Tool

File Replace Options Help

Sequence of replacements:

| Name | Find | Replace with |
|---------------------------------------------------------------|------------------------------------|---------------|
| <input type="checkbox"/> composed - schwa acute with dot | M-X*B*N*o*A*oM-oX*AM-oM-o*AM... | ə̣ |
| <input type="checkbox"/> underlined schwa acute superscrip... | M-C*NM-CM-o*AM-oX*AM-oM-oM... | <u>ə̣</u> |
| <input type="checkbox"/> composed - a acute with dot | M-X*B*N*o*A*oM-oX*AM-oM-o*AM... | ạ́ |
| <input type="checkbox"/> unknown01 - barred i (p20 of c'1... | M-X*B*N*o*A*oM-oG*HM-oM-oC*DM... | &unknown0... |
| <input type="checkbox"/> composed - iota acute | M-X*B*N*o*A*oM-oR*DM-oM-oF*AM... | ị́ |
| <input type="checkbox"/> unknown05 - glottal, line 2102 | M-C*EM-CrM-D*EM-D | &unknown0... |
| <input type="checkbox"/> unknown06 - glottal, line 2141 | M-C*EM-CM-o*HM-oM-D*EM-D | &unknown0... |
| <input type="checkbox"/> composed - i v hook | M-X*B*N*o*A*oM-oD*DM-oM-oC*DM... | &ivhook; |
| <input type="checkbox"/> composed - i acute with dot | M-X*B*N*o*A*oM-oS*AM-oM-oX*AM... | í̇ |
| <input type="checkbox"/> composed - i grave dot | M-X*B*N*o*A*oM-oX*AM-oM-oZ*AM... | î̇ |
| <input type="checkbox"/> composed - u acute with dot | M-X*B*N*o*A*oM-oX*AM-oM-oE*AM... | ú̇ |
| <input type="checkbox"/> composed - i acute with dot | M-X*B*N*o*A*oM-oX*AM-oM-oS*AM... | í̇ |
| <input type="checkbox"/> composed - o acute with dot | M-X*B*N*o*A*oM-oX*AM-oM-o*AM... | ó̇ |
| <input type="checkbox"/> composed - e acute with dot | M-X*B*N*o*A*oM-oX*AM-oM-o*AM... | é̇ |
| <input type="checkbox"/> composed - upsilon acute | M-X*B*N*o*A*oM-oF*AM-oM-o/*AM-o... | ύ̇ |
| <input type="checkbox"/> composed - schwa grave (find in ... | M-X*B*N*o*A*oM-oE*DM-oM-oI*AM-o... | ə̂ |
| <input type="checkbox"/> composed - schwa slash hook | M-X*B*N*o*A*oM-oE*DM-oM-o9*AM-o... | &schwaslas... |
| <input type="checkbox"/> composed - schwa grave | M-X*B*N*o*A*oM-oM-o*LM-oM-o9*AM... | ə̂ |
| <input type="checkbox"/> composed - l ring | M-X*B*K*o*A*oM-oZ*AM-oK*oBM-X | ł̣ |
| <input type="checkbox"/> composed - n vertical line | M-X*B*K*o*A*oM-oM-o*LM-oM-oK*oBM-X | ṇ |
| <input type="checkbox"/> composed - n vertical line | M-X*B*K*o*A*oM-oT*AM-oM-oK*oBM-X | ṇ |
| <input type="checkbox"/> composed - m vertical line | M-X*B*K*o*A*oM-oM-o*LM-oM-oK*oBM-X | ṃ |
| <input type="checkbox"/> composed - m vertical line | M-X*B*K*o*A*oM-oT*AM-oM-oK*oBM-X | ṃ |
| <input type="checkbox"/> modifier barred i | M-C*EM-CM-oC*DM-oM-D*EM-D | ı̄ |
| <input type="checkbox"/> modifier epsilon | M-C*EM-CM-oO*FM-oM-D*EM-D | ε |
| <input type="checkbox"/> modifier a acute | M-C*EM-CM-o*AM-oM-D*EM-D | &modaaacute; |
| <input type="checkbox"/> composed - modifier schwa dot | M-C*EM-CM-oI*AM-oM-D*EM-D | &modschw... |
| <input type="checkbox"/> modifier schwa | M-C*EM-CM-o9*AM-oM-D*EM-D | ə̣ |
| <input type="checkbox"/> composed - modifier i acute | M-C*EM-CM-oS*AM-oM-D*EM-D | &modiacute; |
| <input type="checkbox"/> modifier upsilon | M-C*EM-CM-o/*AM-oM-D*EM-D | υ̇ |
| <input type="checkbox"/> subscript 2 | M-C*FM-C2M-D*FM-D | ₂ |
| <input type="checkbox"/> subscript 1 | M-C*FM-C1M-D*FM-D | ₁ |
| <input type="checkbox"/> modifier y | M-C*EM-CyM-D*EM-D | Ƴ |
| <input type="checkbox"/> modifier u | M-C*EM-CuM-D*EM-D | Ƶ |

Do Replacements ☐ Checked items only

Source text: 010 101

.sf M-)M-o1*AM-oM-oE*AM-o
df reflexive

il M-oS*DM-oM-o+*AM-oM-o1*AM
-oM-o9*AM-o M-oD*GM-oP*M-
o+*AM-oq=M-o1*AM-oM-o9*AM-onM
)M-o1*AM-oM-oE*AM-o
df she is cooking
k S2.106,105

il M-oD*GM-oP*M-X*B*N*o*A*oM-
oE*DM-oM-o9*AM-oN*oBM-XIKM-)
M-o9*AM-onM-o9*AM-onM-)M-
o1*AM-oM-oE*AM-o
df he turned around
k S2a.123; Y6.387

C:\...AFIX3HMG.WRK.txt

Output text: &# &#

.sf -čút
df reflexive

il ʔéčə √p'éq=čan-čút
df she is cooking
k S2.106,105

il √p'əlk-əman-čút
df he turned around
k S2a.123; Y6.387

il √čškw-əman-čút
df he is pulling
k S2a.125

il √málx-əʔan-čút
df he told a lie
k S2a.127

C:\...My Documents\linguistics\ewa\gregs_work\moses\moses-characters.seq.xml

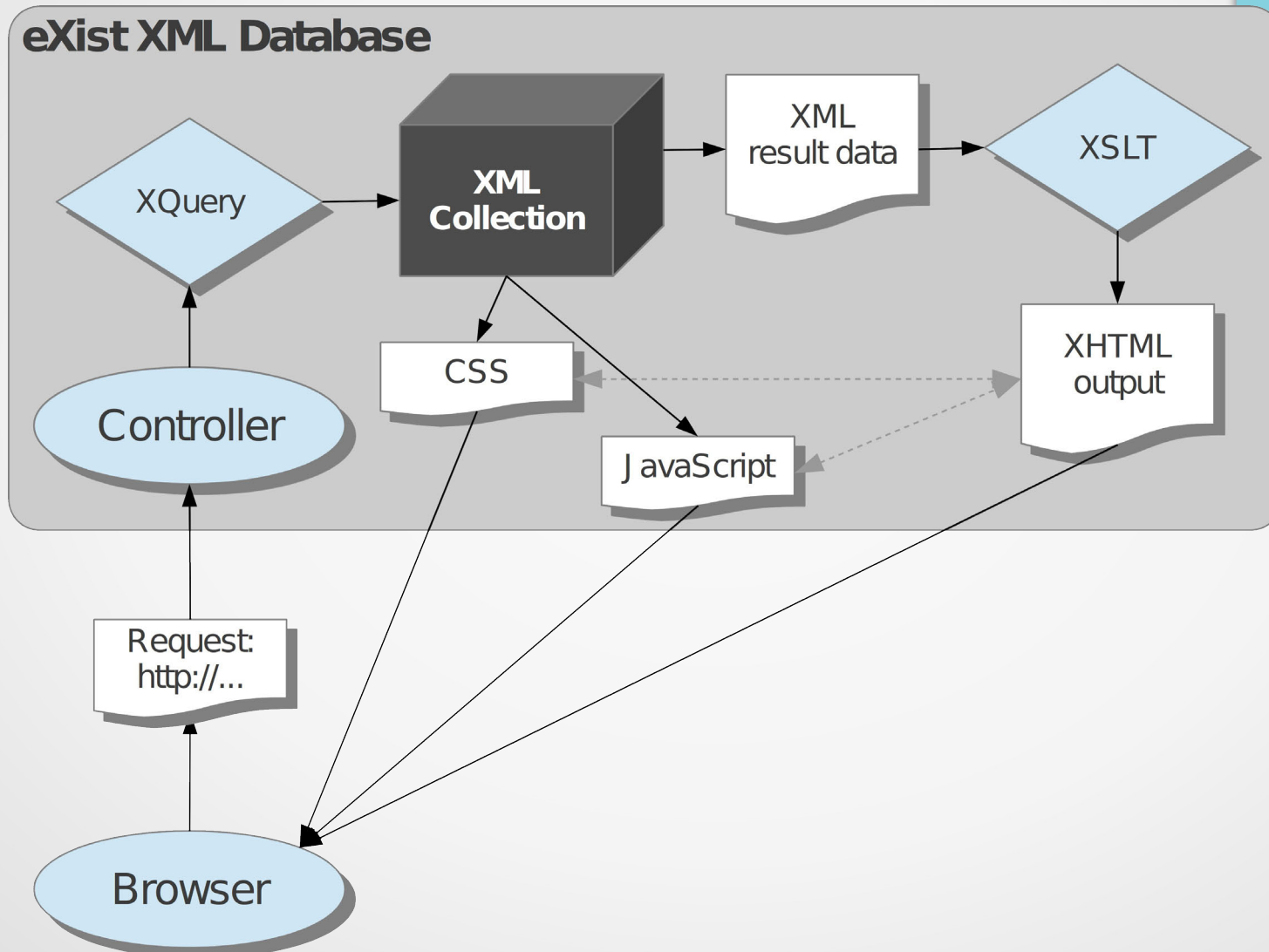
To XML...

```
<ENTRY level="001" id="√cah">
  <rt>√cah</rt>
  <ENTRY level="002" id="">
    <ls></ls>
    <infl mode="1">relational</infl>
    <mn mode="11">√cah=cí-mən</mn>
    <g mode="11">I *encourage|d someone</g>
    <k mode="11">W11.58</k>
    <var mode="11">√çah=čé-mən</var>
    <g mode="11">I *remind|ed him</g>
    <k mode="11">JM3.119.5</k>
    <mn mode="12">√cah=cí-mən-c</mn>
    <g mode="12">he *encourage|d me</g>
    <k mode="12">W11.59</k>
    <mn.m mode="13">√çah=čé-m-əm</mn.m>
    <g mode="13">I *remind|ed him</g>
    <k mode="13">Y41.125</k>
  </ENTRY>
```

To TEI...

```
<entry xml:id="cahcimn">
  <form>
    <pron>
      <seg type="p" subtype="i">cahcímn</seg>
      <bibl corresp="psn:ECH">ECH</bibl>
      <seg type="n">cahcímən</seg>
      <bibl corresp="psn:W">W11.58</bibl>
    </pron>
    <hyph>√<m corresp="m:cah">cah</m>=<m corresp="m:cin">cí</m>-<m corresp="m:min m:t-TR m:Ø-OBJ m:n-SUBJ">mn</m>
    </hyph>
  </form>
  <sense>
    <def>
      <seg>I <gloss>encourage</gloss>d someone</seg>
      <bibl corresp="psn:W">W11.58</bibl>
    </def>
  </sense>
  <xr>See <ref target="m:çahcimn">çahcimn</ref>, and <ref target="m:n?iŋWqn">n?iŋwqn</ref>.</xr>
</entry>
```

The Web application



The printouts

